1RMA: Re-envisioning Remote Memory Access for Multi-tenant Datacenters

Arjun Singhvi, Aditya Akella, Dan Gibson, Thomas F. Wenisch, Monica Wong-Chan, Sean Clark, Milo M.K. Martin, Moray McLaren, Prashant Chandra, Rob Cauble, Hassan M. G. Wassel, Behnam Montazeri, Simon L. Sabato, Joel Scherpelz, Amin Vahdat





Problem Statement

RDMA is an attractive option for modern datacenter applications due to its low latency and high throughput benefits

- Pilaf(atc13), HERD(sigcomm14), FaSST(osdi16), FaRM(nsdi18), PolarFS(vldb18)

Operationalizing RDMA in multi-tenant datacenters uncovered deployment issues:

- Connection-orientedness leading to poor scalability, ordering and failure semantics
- Not amenable to support key security management operations (e.g., rotating encryption keys)
- Not amenable to rapid iteration due to key algorithms (e.g., congestion control) baked in hardware



1RMA: Remote Memory Access for Multi-tenant Datacenters

1RMA NIC Hardware

- Connection-free and fixed-function hardware
- Line-rate crypto for unified access control, integrity and privacy with first class support for encryption key rotation
- Aids software by providing delay measurements and fast failure notifications

1RMA Software

- Implements congestion control and op management (e.g., failure recovery)

RDMA allows direct access to the memory of a remote machine resulting in low latency and high throughput





DRAM















|| Connections State || > RNIC SRAM





Critical applications can be stranded due to connection exhaustion

Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to **induced ordering**

Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering

- RDMA requires FIFO execution of ops (same type) within a single connection



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering

- RDMA requires FIFO execution of ops (same type) within a single connection



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering

- RDMA requires FIFO execution of ops (same type) within a single connection



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering

- RDMA requires FIFO execution of ops (same type) within a single connection



Address connection scalability and performance issues via **connection sharing**

- Multiplexing independent workloads on the same connection

Connection sharing leads to induced ordering

- RDMA requires FIFO execution of ops (same type) within a single connection



RDMA ties access control to connections through protection domains

RDMA does not provide any ready means to manage encryption keys

RDMA ties access control to connections through protection domains

RDMA does not provide any ready means to manage encryption keys



RDMA ties access control to connections through protection domains

RDMA does not provide any ready means to manage encryption keys



RDMA ties access control to connections through protection domains

RDMA does not provide any ready means to manage encryption keys



RDMA Challenge #3: Rigid Congestion Control and Loss Recovery

RoCE enables RDMA over Ethernet fabric and uses priority flow control (PFC) to provide a near-lossless network

PFC is untenable in commercial datacenters due to its operational hazards

- Practical DCB for Improved Datacenter Network(infocom14)
- RDMA over Commodity Ethernet at Scale(sigcomm16)

Recent proposals reduce reliance on PFC

- Bake congestion response in hardware
 - DCQCN(sigcomm15), HPCC(sigcomm19)
- Bake sophisticated loss recovery in hardware
 - IRN_(sigcomm18) : selective retransmission instead of go-back-N

Not amenable to changes post-deployment

RDMA Challenges: Root Causes

RDMA ill-suited due to its two basic design attributes

Connection Complex policies Orientedness baked into hardware

1RMA Approach

Judicious division of labor between hardware and software leading to a simple and fixed-function 1RMA NIC aided by 1RMA software

Fixed-function NIC hardware with explicitly allocated resources

- Connection-free independent ops
- Explicitly-finite hardware resource pools
- Solicitation

Connection-free security protocol with management ops

Software-driven, hardware-assisted congestion control

1RMA NIC acts on **fixed-sized ops** and treats them **independently**

- Provides fail-fast behaviour

1RMA NIC acts on fixed-sized ops and treats them independently

- Provides fail-fast behaviour: NIC ensures op completion within a fixed time



1RMA NIC acts on fixed-sized ops and treats them independently

- Provides **fail-fast behaviour**: NIC ensures op completion within a fixed time



1RMA NIC acts on fixed-sized ops and treats them independently

 Provides fail-fast behaviour: NIC ensures op completion within a fixed time; otherwise delivers fast and precise op failure notifications to software Op from Precise failure



1RMA NIC acts on **fixed-sized ops** and treats them **independently**

- Provides **fail-fast behaviour**: NIC ensures op completion within a fixed time; otherwise delivers fast and precise op failure notifications to software

1RMA NIC leaves retry, ordering, congestion control and segmentation to software

1RMA NIC acts on fixed-sized ops and treats them independently

- Provides **fail-fast behaviour**: NIC ensures op completion within a fixed time; otherwise delivers fast and precise op failure notifications to software

1RMA NIC leaves retry, ordering, congestion control and segmentation to software

RDMA

1RMA



1RMA NIC acts on fixed-sized ops and treats them independently

- Provides **fail-fast behaviour**: NIC ensures op completion within a fixed time; otherwise delivers fast and precise op failure notifications to software

1RMA NIC leaves retry, ordering, congestion control and segmentation to software

RDMA

1RMA



1RMA NIC state does not grow with endpoint pairs

1RMA Key Idea #2: Explicitly-Finite NIC Resource Pools
1RMA Key Idea #2: Explicitly-Finite NIC Resource Pools



Fixed-size Registered Region Table (RRT) and Command Slot Table (CST)

1RMA Key Idea #2: Explicitly-Finite NIC Resource Pools



Fixed-size Registered Region Table (RRT) and Command Slot Table (CST) **Simple Hardware**

Predictable Performance

1RMA Key Idea #2: Explicitly-Finite NIC Resource Pools



Fixed-size Registered Region Table (RRT) and Command Slot Table (CST) **Simple Hardware**

Predictable Performance

Resource allocation based on business priorities



All data transfers in 1RMA are solicited



All data transfers in 1RMA are solicited



All data transfers in 1RMA are solicited



All data transfers in 1RMA are solicited



All data transfers in 1RMA are solicited









































Limits the severity of sudden incasts by bounding the number of inbound bytes





All transfers are **encrypted and signed** in hardware using a novel **connectionfree protocol**



All transfers are encrypted and signed in hardware using a novel connectionfree protocol



All transfers are **encrypted and signed** in hardware using a novel **connectionfree protocol**



All transfers are **encrypted** and signed in hardware using a novel connectionfree protocol



Remote Key Rotation Service

1RMA NIC

CST Outgoing Op FIFOs Incoming Op FIFOs


1RMA Key Idea #4: Connectionless Security with Security Ops



Enables key rotation

- Without placing trust on local software stack
- Without involving local CPU

1RMA offers a Rekey op that provides support for encryption key rotation

1RMA Key Idea #4: Connectionless Security with Security Ops



Enables key rotation

- Without placing trust on local software stack
- Without involving local CPU

Only ops corresponding to the affected memory region fail

1RMA offers a Rekey op that provides support for encryption key rotation

1RMA Key Idea #4: Connectionless Security with Security Ops



Enables key rotation

- Without placing trust on local software stack
- Without involving local CPU

Only ops corresponding to the affected memory region fail

Unavailability period can be reduced to a single RTT

1RMA offers a Rekey op that provides support for encryption key rotation













Uses delay as a congestion signal

Uses delay as a congestion signal

Reacts separately to local and remote congestion

Uses delay as a congestion signal

Reacts separately to local and remote congestion

- Uses a CWND for each remote destination and direction

Uses delay as a congestion signal

Reacts separately to local and remote congestion

- Uses a CWND for each remote destination and direction
- Uses a single CWND for local congestion

Uses delay as a congestion signal

Reacts separately to local and remote congestion

- Uses a CWND for each remote destination and direction
- Uses a single CWND for local congestion

Op rate assigned using the more restrictive CWND


```
Algorithm 1: 1RMA CC reaction to local congestion.
 Input: issue_delay, RTT
 Output: cwnd local
 On Successful Op Completion
    cwnd_local_old \leftarrow cwnd_local
    if issue_delay < TARGET_DELAY then
      ▶ Additive Increase (AI)
      if cwnd local \ge 1 then
         cwnd\_local \leftarrow cwnd\_local + \frac{AI}{cwnd\_local} \triangleright AI = 0.25
      else
     | cwnd local \leftarrow cwnd local + AI
      cwnd\_local \leftarrow min(cwnd\_local, CWND\_MAX)
    else
      ▶ Multiplicative Decrease (MD)
      if no decrease in the last RTT time then
         delta \leftarrow issue\_delay - TARGET\_DELAY;
         cwnd\_local \leftarrow
                                             \triangleright MD = 0.5, \beta = 0.8
          max(1 - \beta \cdot (\frac{delta}{issue\_delay}), MD) \cdot cwnd\_local;
         cwnd\_local \leftarrow max(cwnd\_local, CWND\_MIN);
```

```
Algorithm 1: 1RMA CC reaction to local congestion.
 Input: issue_delay, RTT
 Output: cwnd local
 On Successful Op Completion
    cwnd_local_old \leftarrow cwnd_local
    if issue_delay < TARGET_DELAY then
      ▶ Additive Increase (AI)
      if cwnd local \ge 1 then
         cwnd\_local \leftarrow cwnd\_local + \frac{AI}{cwnd\_local} \triangleright AI = 0.25
      else
     | cwnd local \leftarrow cwnd local + AI
      cwnd\_local \leftarrow min(cwnd\_local, CWND\_MAX)
    else
      ▶ Multiplicative Decrease (MD)
      if no decrease in the last RTT time then
         delta \leftarrow issue\_delay - TARGET\_DELAY;
         cwnd\_local \leftarrow
                                             \triangleright MD = 0.5, \beta = 0.8
          max(1 - \beta \cdot (\frac{delta}{issue \ delau}), MD) \cdot cwnd_local;
         cwnd\_local \leftarrow max(cwnd\_local, CWND\_MIN);
 On Dispatch Timeout
    if no decrease happened in the last RTT time then
      cwnd \ local \leftarrow
                                       \triangleright TIMEOUT DECR = 0.1
```

```
TIMEOUT\_DECR \cdot cwnd\_local;
cwnd\_local \leftarrow max(cwnd\_local, CWND\_MIN)
```

Evaluation

How does 1RMA perform relative to state-of-the-art?

Workload: 40-node uniform random traffic pattern Baselines: Pony Express(sosp19) and RDMA

1RMA experiences the least latency slowdown for comparable offered load

Does 1RMA provide isolation and ensure prioritization?

Point-to-point workload consists of small 64B foreground transfers and one outstanding large background transfer

RDMA introduces delays proportional to the background RMA size while 1RMA experiences minimal slow down

What impact does 1RMA's Rekey Op have on availability?

Point-to-point workload consists of small reads and server does encryption key rotation

1RMA experiences minimal impact of availability during key rotation

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Client initiates long-running transfers of 4MB reads with one server in the beginning and another server at around 400 us

No separate reaction to local congestion

Separate reaction to local congestion

Reacting separately to local congestion converges 20x faster

Separate reaction to local congestion

No separate reaction to local congestion Submission Rate (Gbps) Server 1 100Server 2 75 50 25 0 200 400 600 800 1000 1200 1400 0 Time(us)

Summary

Existing RDMA technologies ill-suited to multi-tenant datacenters

- Connection-orientedness
- Complex policies in hardware

1RMA re-envisions RMA for multi-tenant datacenters

- 1RMA NIC hardware
 - Connection-free and fixed-function with explicitly allocated resources
 - Aids software by provided delay measurements and fast failure notifications
- 1RMA software
 - Implements congestion control and op management (e.g., failure recovery)

Thank you!

1RMA: Re-envisioning Remote Memory Access for Multi-tenant Datacenters

Arjun Singhvi, Aditya Akella, Dan Gibson, Thomas F. Wenisch, Monica Wong-Chan, Sean Clark, Milo M.K. Martin, Moray McLaren, Prashant Chandra, Rob Cauble, Hassan M. G. Wassel, Behnam Montazeri, Simon L. Sabato, Joel Scherpelz, Amin Vahdat

